

```

<!DOCTYPE HTML PUBLIC "-//W3C//DTD HTML 4.01//EN"
"http://www.w3.org/TR/html4/strict.dtd">
<html>
<head>
<title>IEEE 754 and me - an experiment in "anonymizing" a PDF file</title>
<meta http-equiv="Content-Type" content="text/html; charset=ISO-8859-1">
</head>
<style type="text/css"><!--
i, pre { background: #ccf }
i.white { background: white }
--></style>
<body style="max-width:450px; text-align:justify">

<h1>IEEE 754 and me - an experiment in <i class="white">anonymizing</i>
a PDF file</h1>

<p>(I like to create words.)
<p>In this page I explain a hack I did with the PDF file that
contains the IEEE 754 standard.
<p>I don't have access to the IEEE website, so I normally would have
to pay to read this file. Except I know some people who have access.
(I mean: access to some publications. Obviously if I enter the IEEE
address in any browser I'll get some data back in my face.)
So I just asked to them if they could get it for me.
<p>So they did. And they gave it to me.
<p>Except...
<p>Except if I want to redistribute the file, I can't.
<p>Well I can, obviously. It's just a file. I can put it on a website and
there we are, it's shared.
<p>But my problem is that the file contains some data from the people who
gave it to me. So <i class="white">they</i> would be in trouble if IEEE
finds the file. They just look at it and "oh, but these are these people
who distributed the file to the world. Let's sue them!"
<p>I am close to these people. I just can't do that.
<p>Except...
<p>Except, you know, a file is a file. You can edit a file. You can
modify it. You can extract stuff and put new stuff inside. You can
do whatever you want with the file.
<p>So you can remove the stuff that identify those who gave this
particular PDF file to me.
<p>That's what I did.
<p><a href="ieee754.tar.gz">Here</a> is some code to prove my claim.

<p>Let me give some more details, will you my dear? It might be useful
to young hackers out there. Some kind of lesson of hackery. I use
the C language. I work under UNIX.

<p>(Ah no. I don't distribute the clean IEEE 754 file.)
<br>(I AM NOT A CRIMINAL!)
<br>(Well...)

<h2>What is a PDF file?</h2>

<p>A PDF file is a file that conforms to a given structure,
namely the PDF one.

```

<p>A description of it is freely available from Adobe.
See here.
(Adobe, like many big private corporations and many institutions, changes
its internet website's structure from time to time. You might need to
dig a bit to get the specifications. Alternatively you can look for
"PDF specifications" with your favorite search engine.)

<p>Since the PDF file I want to hack has a version of 1.4, I went
there
to download
that.
The current version, as of writing, is 1.7. Don't ask me why people
like to change stuff when everything works fine. It's called
"business" and it's something I'm completely unable to understand since
it is so irrational.

<p>So by reading this (big) documentation, you understand the structure
of a PDF.

<p>A PDF is made of several parts.

<p>First part: the header. In our case, the first line of the file contains
the version (<i>%PDF-1.4</i>) and the second contains a comment with
byte greater than 127. These bytes are here, as prescribed by Adobe,
to inform software out there that the file contains binary data, ie.
bytes greater than 127. Why does it matter? Because the internet is an
old beast and in the past (and maybe still now) some links or protocols
or both send only bytes smaller than 127 in some cases. Putting high
bytes early in the file forces them to switch to a full 8 bits mode,
thus accepting bytes from 0 to 255 (this might not be correct, my
memory is a bit, you know, well, bad).

<p>Anyway, what matters to us here is that these bytes are normally
chosen at random. But how do we know? So the file we are going to
produce won't contain these numbers, but some others, just in
case IEEE attach some meaning to them.

<p>After the header you (generally) have objects, a lot of them.
In our case we have more than 700 objects. What does an object
look like? Well, something like the following.

```
<p>  
<pre>  
2 0 obj<&lt;&lt;/Type/Pages/Kids[3 0 R 4 0 R 5 0 R 6 0 R 7 0 R]/Count  
70&gt;&gt;&gt;  
endobj  
</pre>
```

<p>You find a first number, which is the object's ID, starting from 1
and counting up, then a 0 (generally, can be used for stuff we don't
care here) (all the objects in the IEEE file have a 0 here). Then
you have <i>obj</i> then some stuff. Then, on the same or another line,
you have <i>endobj</i> to, guess what? end the definition of an object.

<p>Some objects may be compressed, which complicates slightly the
extraction, but that's not very hard either.

<p>Note that an objet's definition does not necessarily has to fit in one line. In our case IEEE put almost everything on the first line, which eases our life, but it's not mandatory.

<p>After all the objects, you find the trailer, which contains something called <i>xref</i>. This stuff lists all the objects with their starting offset in the file. Then you find some more stuff useful for programs so you must keep them. One of the things in there is called <i>/ID</i> and this one looks strange so we will remove it. Maybe IEEE tracks the file with the help of this string.

<p>And that's it for a PDF file. It's nothing particularly tricky to understand and manipulate. You must record where in the file each object resides. And if you want to modify an object it changes the position of all the following objects. And this is a very bad design if I dare. But my bet is Adobe doesn't care I don't like their design.

<h2>Extracting objects</h2>

<p>So we want to modify this IEEE file, to remove stuff put in there to identify who downloaded it, to "anonymize" the file.

<p>How do we do that?

<p>Well, after some thinking the conclusion pops up, obvious. We have to extract all the objects, modify those we have to and glue back everything together when we are done.

<p>We can't modify the file directly.

<p>Why not? I mean, we just could overwrite the data we don't want in and that's it, we are done.

<p>But it's not that simple.

<p>First, the size of a file is an information per se. So keeping it might be bad. We will <i class="white">remove</i> objects (or more specifically objects' content since removing an object is a pain; you must then remove all references to it and this requires a lot of parsing).

<p>Then some objects are compressed. Uncompressing them and storing them back changes their size.

<p>(Well, after a bit of thinking, this does not appear to be that important. Emptying objects reduces their sizes. And an uncompressed object is described with less characters. So we could simply put an empty version of the object and the overwrite the remaining space with spaces.)

<p>So we are up to extract object. And this is not very easy at first.

<p>And here comes a bit of hacking.

<p>The IEEE file is very friendly. Objects are all stored the same way. At the start of a line, you have the object's ID, a 0 and <i>obj</i>. And they start right after the header. So, just bypass the header, and read objects, up to a point where what you read is not an object anymore.

<p>And since objects start with a simple line, just read one line. Check for the presence of <i>obj</i> in it. If it is present you have an object and keep going. If not you are done with the objects and you are up with the trailer.

<p>And you know what? It's even simpler than that. After you smoke the header, read just one byte. If it's <i>x</i> you are in the trailer, if not it's an object.

<p>How do I know that? By looking inside the file. A <i>cat file | less</i> is very useful. If you need to look at the binary data, I suggest using <i>od</i> like <i>od -tx1 -ta file | less</i> (do a <i>man od</i> to get more options).

<p>So I looked at the file and just saw that to discriminate between the start of an object and the trailer I just needed to look for an <i>x</i>.

<p>I normally would have to scroll a lot to get to the trailer when I examine the file. Fortunately, there is <i>tail</i> a program that only shows the last lines of a file. To see the last 400 lines of a file you do <i>tail -n 400 file | less</i>. This saves a lot of scrolling time.

<p>Here comes some code I wrote based on all that.

<p>First, the main function.

```
<p>
<pre>
int main(void)
{
    header();
    while (obj()) dump_obj();
    trailer();
    return 0;
}
</pre>
```

<p>See, it's very short. It first deals with the header, then gets the objects, then the trailer.

<p>By looking at the file, I know the header is only made of two lines, so here is <i>header</i>.

```
<p>
<pre>
void header(void)
{
```

```

FILE *f = fopen("header.obj", "wb");

ERR(f, "header");

read_line();
dump_buffer(f);
read_line();
dump_buffer(f);

if (fclose(f)) ERR(0, "header");
}
</pre>

```

<p>Very clear too.

<p>For <i>obj</i> things are a bit more complicated.

```

<p>
<pre>
int obj(void)
{
    int c;
    c = getchar(); eof(c);
    if (c == 'x') { ungetc(c, stdin); return 0; }
    no++;
    size = 0;
    put_char(c);
    while (1) {
        c = getchar(); eof(c); put_char(c);
        if (c == '\n') if (endobj()) break;
    }
    return 1;
}
</pre>

```

<p>As you can see, we first check for an <i>x</i>. Then we read everything until we see <i>endobj</i>. Once again, by inspecting the file, I saw that <i>endobj</i> only appears at the beginning of a line, so I look for it only when I received an "end of line" character, <i>\n</i>.

<p>What does <i>endobj</i> do? It just reads seven characters. (Seven and not six because I also saw that after <i>endobj</i> there is an "end of line" character which I include, just to get the correct start for the next object, or the <i>x</i> starting the trailer). If at some point it doesn't get the expected byte, it returns <i>0</i>, meaning we don't read <i>endobj</i> yet. Otherwise it returns <i>1</i>.

```

<p>
<pre>
int endobj(void)
{
    int c;
    c = getchar(); eof(c); put_char(c); if (c!='e') return 0;
    c = getchar(); eof(c); put_char(c); if (c!='n') return 0;
}

```

```
c = getchar(); eof(c); put_char(c); if (c!='d') return 0;
c = getchar(); eof(c); put_char(c); if (c!='o') return 0;
c = getchar(); eof(c); put_char(c); if (c!='b') return 0;
c = getchar(); eof(c); put_char(c); if (c!='j') return 0;
c = getchar(); eof(c); put_char(c); if (c!='\n') return 0;
return 1;
}
</pre>
```

<p>One important point here: if <i>endobj</i> appears in an object we are doomed. I played the game, because it makes my life much easier. And I won! No <i>endobj</i> in weird places in this file. Otherwise more complicated parsing would have been required. And you may not know it since you don't know me but I'm a very lazy person.

<p>So that's it for the objects.

<p>I pass the trailer's extraction. We simply read everything past the last object and put it in a file.

<h2>Uncompressing compressed objects</h2>

<p>Once we extracted the objects we have to inspect their content, one by one, to decide what to do (keep the object, empty it or modify it).

<p>And there comes a problem. Some objects are compressed. We see something like the following when using <i>less</i>.

```
<p>
<pre>
331 0 obj<&lt;&/Filter/FlateDecode/Length 10&gt;&gt;stream
x<&lt;9C&gt;+&lt;E4&gt;^B^@^&lt;EE&gt;^@|
endstream
endobj
</pre>
```

<p>Obviously, we have some more work to do.

<p>And there comes something very important for a hacker. You just cannot face a problem like this one and solve it out of nothing, as by magic. You need something. You need enough knowledge to get the right intuition at how to attack the problem.

<p>You need, in a word, <i class="white">culture</i>.

<p>So here you go read the PDF specifications to learn about this <i>/Filter/FlateDecode</i> stuff. You read, page 46, that this encoding method uses <i>zlib/deflate</i>.

<p>And there comes your culture. Because, what to do of this knowledge? What if you don't know what this <i>zlib</i> stuff is?

<p>By reading others' code, by browsing a lot of software on your system and in the internet, by compiling by hand, by carefully looking at what happens then, and also, of course, by using software libraries,

you build a knowledge base in yourself. You become sort of an expert.

<p>And there, in the case of <i>xlib/deflate</i>, I know what it is, and I know how to handle that. It means that the data has been compressed and that I can uncompress it using the popular <i>xlib</i> library, which everyone uses or almost. And it's almost certainly installed on your computer. Dozens and dozens of software you use on a daily basis use it. And you can use it too.

<p>So you read <i>/usr/include/zlib.h</i> (which is a bit of a pain, but you can manage it) and after a few trials and errors, you quickly write some code to uncompress a zlib's stream.

<p>And you end up with something like the following, extracted from <i>deflate.c</i>. (Yes, I reversed inflate/deflate in my brain. Normally <i class="white">deflate</i> is to decrease the size of a file, but I wrongly thought it meant <i class="white">decompress</i>. Poor french native speaker of me...)

<p>

<pre>

```
    strm.next_in = (unsigned char *)buf;
    strm.avail_in = size;
    strm.next_out = (unsigned char *)out;
    strm.avail_out = 4000000;
    strm.zalloc = 0;
    strm.zfree = 0;

    if (inflateInit(&strm) != Z_OK)
        { fprintf(stderr, "inflateInit failed.\n"); return 7; }
    if (inflate(&strm, Z_FINISH) != Z_STREAM_END)
        { fprintf(stderr, "inflate failed.\n"); return 8; }
    if (inflateEnd(&strm) != Z_OK)
        { fprintf(stderr, "inflateEnd failed.\n"); return 9; }
```

</pre>

<p>It looks like magic? It's not. We setup a <i>strm</i> structure, then we call <i>inflateInit</i>, <i>inflate</i> and <i>inflateEnd</i> in order and check that the return value is correct.

<p>I had troubles at first for <i>inflate</i>, not checking against <i>Z_STREAM_END</i> but <i>Z_OK</i>. I read <i>/usr/include/zlib.h</i> too fast.

<p>It might happen to you as well. You may read documentation too fast, or understand it wrongly, leading to some time trying to understand what's wrong. It may become very frustrating after some time, especially when you don't hack just for yourself, with no time pressure, but when you work in a company, with a tight schedule. Key here is to be patient. In most of the cases, when using some new library, the problem comes from you. In some cases (I had the problem once a few years ago with <i>ptrace</i>; the manpage was wrong) the error is in the documentation. In some other cases it's a bug in the software. How to know

where the problem comes from so that you spend more time solving it and less time tracking it? It's hard to answer. There again culture plays a role. When you start your hackery's journey, you know nothing. So you spend a lot of time doing useless things. After a while (several years) you start to understand what's going on and you focus on the important stuff.

<p>In my case, I more or less applied the following reasoning to dig the problem. Since <i>zlib</i> is a well know library, since PDF files are also well established, it is unlikely that one or both of these beasts have a problem. The problem very certainly comes from me. What did I do wrong then? I read <i>/usr/include/zlib.h</i>, maybe a bit too quickly. So maybe I should read it again, more carefully. So I did. And I then found that using <i>Z_OK</i> was inappropriate in my case. I had to use <i>Z_STREAM_END</i>.

<p>That's how you look for a bug in your software. First, try to remember how you wrote it. Then value each step with a "confidence" number. You obtain a list of actions, the first being the most probably wrong and the last being practically 100% correct. (Normally you only have very few "most probably wrong" actions at a given point of your coding. Organizing your work is part of the journey too. It's a skill to get.) Then you review actions in order, looking carefully where you originally just looked quickly at some stuff.

<p>You need to be humble. Most of the time <i class="white">you</i> are the cause of the error. <i class="white">You</i> didn't understand something clearly enough. <i class="white">You</i> typed something wrong on your keyboard. <i class="white">You</i> called some functions in the wrong way, in the wrong order, with wrong arguments. <i class="white">You</i> didn't check return values correctly and implemented an appropriate response (like ending the program or displaying a meaningful message on the terminal).

<p>But you must also trust your judgment. Sometimes you review all your possible mistakes and nothing comes out of it. So the problem is from somewhere else. It might be a library, the compiler, the language you use, the operating system, the hardware, anything. How do you know that you are not guilty and that the problem is from another part of the system? And which part? How do you decide you digged enough in one direction that you are almost confident the problem is not here?

<p>There is no clear answer. One important thing to do when facing a problem is trying to decrease its size, extract it. You have a bug in your 1,000 lines program? Try to write a very small one that also exhibits the problem. The key point is the ability to extract a "model" on which you can think with full confidence. Something like "okay, I extracted a small version of my program; I am now 100% sure that if

I don't find a problem with this little program then my bigger program is not faulty but the sytem, somewhere, is. Then I'll try to narrow it to the correct part of the system, eliminating, one by one, every possible cause, up to the point I am left alone with the guilty beast. I will then shoot it in the head and I will be left with a perfectly correct program, solving my problem as it is intended."

<p>And that takes a lot of time, a lot of trials and errors, a good analytical brain and a huge dose of patience and humility.

<p>And here comes the <i>main</i> function of <i>deflate.c</i>.

```
<p>
<pre>
int main(void)
{
    char *t;
    char *l1, *l2;
    int l;
    int c;
    z_stream strm;

    read_line();
    put_char(0);

    if (!strstr(buf, "obj&lt;&lt;"))
        { fprintf(stderr, "WARNING: no obj found.\n"); return 1; }

    if (!strstr(buf, "&gt;&gt;stream"))
        { fprintf(stderr, "no stream found.\n"); return 2; }

    if (!strstr(buf, "/Filter/FlateDecode"))
        { fprintf(stderr, "no filter flatdecode found.\n"); return 3; }

    if (!(t=strstr(buf, "/Length "))
        { fprintf(stderr, "no length found.\n"); return 4; }

    if (sscanf(t, "/Length %d", &l) != 1 || l &lt; 0)
        { fprintf(stderr, "bad length found.\n"); return 5; }

    fprintf(stderr, "INFO length %d\n", l);

    l1 = strdup(buf); ERR(l1,"strdup");
    l2 = malloc(size); ERR(l2,"malloc");

    size = 0;
    for (; l; l--) {
        c = getchar(); eof(c); put_char(c);
    }

    strm.next_in = (unsigned char *)buf;
    strm.avail_in = size;
    strm.next_out = (unsigned char *)out;
    strm.avail_out = 4000000;

```

```

strm.zalloc = 0;
strm.zfree = 0;

if (inflateInit(&strm) != Z_OK)
    { fprintf(stderr, "inflateInit failed.\n"); return 7; }
if (inflate(&strm, Z_FINISH) != Z_STREAM_END)
    { fprintf(stderr, "inflate failed.\n"); return 8; }
if (inflateEnd(&strm) != Z_OK)
    { fprintf(stderr, "inflateEnd failed.\n"); return 9; }

fprintf(stderr, "INFO output bytes %ld\n", strm.total_out);

new_first_line(l1, l2, strm.total_out);

printf("%s", l2);

fwrite(out, strm.total_out, 1, stdout);

while (1) {
    c = getchar(); if (c == EOF) break;
    printf("%c", c);
}

fflush(stdout);

return 0;
}
</pre>

```

<p>We look at the first line in the file. By inspecting the extracted objects with <i>less</i> we found out that the first line contains all what is necessary to decide if the object is compressed or not. That's what we do, that's what all the tests do.

<p>After that, if we indeed have a compressed file, we get the length of the compressed data, read this data, decompress it, and put everything back in a new, uncompressed, object. I don't show you <i>new_first_line</i>, it just removes <i>/Filter/FlateDecode</i> and changes the <i>/Length</i> stuff.

<p>One interesting point here. The program will produce an empty output (and some messages on <i>stderr</i>). And here is how I run it on all the produced objects.

```

<p>
<pre>
for i in ????.obj; do ./deflate &lt; $i &gt; out/$i; done
</pre>

```

<p>So in <i>out</i> some files will be empty. We can use that information to overwrite empty files (which thus are not compressed) with the current version. Here is how we do that.

```
<p>
<pre>
cd out; find -empty -exec cp ../{} . \;
</pre>
```

<p>This line (and maybe the previous one) may look a bit magic.

<p>There again we come back to *culture*. These commands (*for*, *find*, ...) and their arguments are "common knowledge" you acquire over time, by reading documentation, trying the command by yourself, and looking at how other people use them. There is no magic, just culture.

<p>Once it's done you are with objects in *out*, that you finally can fully inspect and modify.

<h2>Modifying objects</h2>

<p>In writing this webpage, I chose a linear order to present steps involved in the analysis and modification of a PDF file. The reality is that after I extracted the objects, and even before uncompressing them, I wrote *glue.c* that takes them all back and glues them together in a new PDF file. So maybe I should have put the "gluing objects" before this section. I don't know. When you hack you often write stuff "out of order" so as to have a few problems to solve at once. By writing *glue.c* early, just after the extraction, I could check that the extraction was correct, that the resulting PDF file obtained by collecting all the objects was similar to the original PDF. In fact, I needed to know if the extraction was correct. So obviously, the way to test that is to put them all back together and see what happens. And since I would need such a tool later anyway, why not write it right now? So that's why I wrote *glue.c* right after *extract.c*.

<p>But let's look a bit at what needs to be done with our uncompressed objects.

<p>I read them all. All the objects. More than seven hundreds of them. Just to see what was inside, what was needed to be removed to anonymize it.

<p>I quickly found that some objects contained nothing but a text identifying who download the file.

<p>A little surprise was to find several of those objects, all identical. Why not put only one? It's a clear loss of space. Maybe they thought it would be harder to remove. Maybe they just used some software that was suboptimal. It's not a big issue, but it forces us to do more than just editing by hand one file. We

have to come up with an automatic solution.

<p>So here is the magic line I quickly found.

```
<p>
<pre>
cd out; for i in `grep -i limited *|cut -d : -f 1`; do ../empty $i > $i; done
</pre>
```

<p>And here comes <i>empty.c</i>.

```
<p>
<pre>
#include <stdio.h>

int main(int n, char **v)
{
    int i;
    sscanf(v[1], "%d", &i);
    printf("%d 0 obj <</Length 0>>stream\nendstream\nendobj\n", i);
    return 0;
}
</pre>
```

<p>The weird command line does a <i>grep</i> (very common and useful command in the UNIX world) over all the files, looking for "limited" (I first checked that "limited" only appeared in interesting files, and it did, so it's safe to discriminate files based on this string) and getting the file's name out of <i>grep</i>'s answer.

<p><i>grep</i> normally replies something like the following.

```
<p>
<pre>
0025.obj: A B C D limited E F G H
</pre>
```

<p>We have the file's name, double dot, and the line of file that matches what we look for.

<p>We are interested in the file's name, and only that.

<p>The wonderful <i>cut</i> command comes at help here. We ask it to cut the first field <i>-f 1</i> by using the double dot separator <i>-d :</i>. The <i>|</i> (pipe) is a common mechanism under UNIX to plug one command to another. The output of the first command becomes the input of the other. You can chain many commands that way, manipulating data the way you want.

<p>So after <i>grep</i> and <i>cut</i> we have a list of files. For each of those we call empty, that creates an empty object. It's important to pass the file's name to <i>empty</i> for the object's number is written in the file, and that changes, obviously, with each object.

<p>Looking at <i>empty.c</i> we see <i>sscanf(v[1], "%d", &i);</i> that will get the object's number from the file's name.

<p>A last note on the command line. We see <i>for i in `XX` ...</i>, with backquotes. Those backquotes are very useful. They transform the output of the command <i>XX</i> (or list of commands piped together) into arguments in the command line for the <i>for i</i> stuff.

<p>It would have been possible to use <i>xargs</i> here, like in the following.

```
<p>  
<pre>  
grep -i limited *|cut -d : -f 1|xargs ../empty  
</pre>
```

<p>But I am more familiar with the first version (which may be slower, or not even work at all if too much data comes out of <i>XX</i> since command line's arguments have a limited size normally).

<p>So, you have many ways to solve a problem. You don't need to learn all of them. What matters is to know one that works, and one that works fast enough. I could have done everything is C programs, which is my favorite horse. But I know some shell commands, and it's faster to use them than to write a C program from scratch. Maybe someone with a deeper knowledge of the shell and available commands would have done what I did faster. Maybe using other languages, like perl or python, would also speed up the coding of all this.

<p>After some time you stick to your tools. Here, you must be careful. Are your tools good enough for your needs? Shouldn't you learn a few other tools? Where is the limit? When stop learning?

<p>No clear answers there either. It's a matter of taste. And time. If you can solve a problem in one day with your tools and if you want to solve it in less than a week then it's okay. If you need one week of work for something you want to solve in one day you are in trouble. If you have no time constraints, do what you want. Hacking normally involves freedom, no tight schedule, so time should not be a problem. You normally care about beauty, elegance, this kind of stuff. Efficiency is a very subjective matter in the hacking world.

<h3>Removing images and colormaps</h3>

<p>I don't trust images. Neither should you. It's easy to embed hidden data in images.

<p>So I decided to remove them from the file.

<p>Objects containing images are easy to find. There is <i>/Image</i> on the first line (in the IEEE PDF file; remember that this file is very friendly for when you want to parse it; a lot of information is present on the first line).

<p>I found two images. How did I process them?

<p>I just changed <i>/Width</i> and <i>/Height</i>, set them to <i>1</i> and put a <i>0</i> in the stream.

<p>There also was a colormap. I don't trust that either. Maybe (but I seriously doubt it) there is some hidden data there too. So I changed it to be full white, just in case.

<p>Then I found some <i>ID</i> stuff here and there, which I also removed, and a file containing XML data, which I completely removed.

<p>All in all, I tried to remove everything that was not necessary and which I thought would pose an anonymity problem.

<p>Well, the process was not fully done.

<p>I didn't change text's objects. I should have changed the spacing of commands in there. And maybe also their order. I should have changed the number of objects, their order in the file.

<p>I mean, information is everywhere. You can hide a message in many many ways. If you want to do things properly, you have to change everything.

<p>The best in our case would be to transform the PDF file by printing it as images, like <i>png</i> images for example, glue them in a file and distribute that file. Or extract text from them through OCR (optical character recognition), or even text directly from the PDF file, and generate a new PDF based on that.

<p>Yes, to be really anonymous, you can't trust the file you download. You have to extract the data you care about, and be prepared that any other data may contain some information useful to identify you.

<p>In my case, I don't care that much. I seriously doubt they go that far at IEEE to tag a poor file that costs a few dollars and which content is already known by many people, many websites, many sources.

<p>But for images, no. I don't trust images. Period.

<p>But this process could not be automatized. I would have had to do it each time I modify something in my programs.

So I just included the few objects in the sources, and in `transform.sh` I do the following.

```
<p>
<pre>
for i in *.n; do cp $i out/`echo $i|sed -e "s/\.n//"`; done
</pre>
```

Once again, a bit of shell's magic.

This time, what we have to do is to overwrite files in `out/` with hand modified versions. These modified files are named with a `.n` suffix and are in the current directory.

Here the wonderful `sed` command shows up. I love this one. Even if I don't master it at all. So here I want to take a file's name like `0123.obj.n` and make it `out/0123.obj` so that `cp` can do its job (to copy one file to another). So I ask `sed` to look for `.n` and replace it with nothing. I use the backquote trick to get the name back on the command line and `cp` proceeds happily.

Note that I also could have put the modified files in a `mods` directory and do a simple `cp mods/* out/` but my brain, at that particular time, was not in this mood. And the command line above came out quickly, so I wasted no time here. I do this kind of manipulations so often that it's hardwired up there in the brain.

That explains why you sometimes see very convoluted solutions to simple problems. People are used to think in a way, to write code with a given style they're used to, that it ends up as a mess only them can understand.

How to protect against that?

Easy. First, solve your problem. Second, look at your solution with a critical eye. Third, rewrite your solution to be more elegant.

Obviously this all requires time and commitment. In a private company you generally don't do that. You have no time, schedule is tight, your thinking and coding costs a lot. So once you have a solution you stick to it. This is a very bad way to act. But hey, some people made some choices at points in history. Private corporations are one of those choices. They are clearly suboptimal ways to organize work. But they work, somehow, so as long as it costs less to organize work this way and getting suboptimal solutions out of this setup, well, it will keep going. It's very hard to change such a big system.

This is some kind of hacking in itself, one might say. Change the world. On a large scale.

But we get out of programs and computers there,
so let's stay with our simple PDF problem for the
purpose of this webpage.

Recompressing objects

We now have our objects, cleaned up. No identification is possible anymore. (Unless IEEE uses complex schemes to mark its files, which I doubt a lot. They might think a PDF file is so complex no one would ever do the kind of things explained in this webpage.)

It's time to recreate a PDF file from those object.

A first attempt was made, where we simply put them back at this point, without compression whatsoever. The produced file was too big (around 2MB). So a recompression's step was necessary.

How to recompress? Well, straightforward. Just check if the object contains a stream. If yes, then compress the stream. If it turns out that the compressed version is smaller, put it in place.

Note that in the original file streams are compressed even if the compressed size is bigger. It happened for very small streams, but it happens nevertheless. What software they use I don't know, but it does not work very well, if I may.

With no more useless words, here comes the file `inflate.c` that compresses a file if it is compressible.

```
<p>
<pre>
#include <stdio.h>;
#include <stdlib.h>;
#include <zlib.h>;
#include <errno.h>;
#include <string.h>;
#include <ctype.h>;

#define ERR do { printf("%s:%d: err\n", __FUNCTION__, __LINE__); exit(1); }
while (0)

char *in;
int size;
int maxsize;

void put_char(char c)
{
    if (size == maxsize) {
        in = realloc(in, maxsize+=4096);
        if (!in) ERR;
    }
    in[size] = c;
    size++;
}

char out[4000000];
```



```

int main(int n, char **v)
{
    int c;
    FILE *f;
    int l;
    char *s;
    char *t;
    z_stream strm;
    int i;

    if (n != 2) { printf("gimme a file to process\n"); return 1; }

    f = fopen(v[1], "r");
    if (!f) ERR;
    while (1) {
        c = getc(f); if (c == EOF) { if (errno) ERR; break; }
        put_char(c);
    }
    put_char(0);
    fclose(f);
    if (!(t=strstr(in, "&gt;&gt;stream"))){printf("%s has no
stream\n",v[1]);return 0;}
    s = strstr(in, "/Length "); if (!s) ERR;
    if (sscanf(s, "/Length %d", &l) != 1) ERR;
    if (l == 0) {printf("%s has Length = 0, nothing done\n",v[1]);return 0;}
    t+=9;
    if (size &gt; 4000000) ERR;

    strm.next_in = (unsigned char *)t;
    strm.avail_in = l;
    strm.next_out = (unsigned char *)out;
    strm.avail_out = 4000000;
    strm.zalloc = 0;
    strm.zfree = 0;

    if (deflateInit(&strm, 9) != Z_OK)
        { fprintf(stderr, "deflateInit failed.\n"); return 7; }
    if (deflate(&strm, Z_FINISH) != Z_STREAM_END)
        { fprintf(stderr, "deflate failed.\n"); return 8; }
    if (deflateEnd(&strm) != Z_OK)
        { fprintf(stderr, "deflateEnd failed.\n"); return 9; }

    if (strm.total_out &gt; l) {
        printf("compressed size (%ld) is more than uncompressed (%d), do
nothing\n",
            strm.total_out, l);
        return 0;
    }
    f = fopen(v[1], "w"); if (!f) ERR;

    t = in;
    while (t != s) {
        fprintf(f, "%c", *t);
        t++;
    }
}

```

```

    fprintf(f, "/Filter/FlateDecode/Length %ld", strm.total_out);
    while (!isdigit(*t)) t++;
    while (isdigit(*t)) t++;
    while (*t != '\n') {
        fprintf(f, "%c", *t);
        t++;
    }
    fprintf(f, "\n");
    for (i=0; i<&strm.total_out; i++) fprintf(f, "%c", out[i]);
    fprintf(f, "\nendstream\nendobj\n");
    fclose(f);
    return 0;
}
</pre>

```

<p>Note that this time, a bit different from how <i>deflate.c</i> works, we output in the file itself, not on <i>stdout</i>, which eases the life in the shell's world. The simple following command does the trick very well.

```

<p>
<pre>
cd out; for i in ????.obj; do ../inflate $i; done
</pre>

```

<p>I know, I know. My programs <i>deflate.c</i> and <i>inflate.c</i> have wrong names. But now it's done and an important notion in the hacking world is: laziness.

<h2>Gluing objects</h2>

<p>And now for the last action: glue everything back together.

<p>Things are a bit more complex. Or more specifically, verbose. It's not very complex.

<p>Here is the <i>main</i> function.

```

<pre>
int main(void)
{
    get_root_info(&f);
    get_files(&f);
    dump(&f);
    return 0;
}
</pre>

```

<p>We need to get the "root" and "info" stuff from the trailer, to let the new PDF file start at the same position than the original one. What does it mean for a PDF file to "start?" Well, we have a list of object. And one of them is the "root" of everything. This has to be specified in the trailer. For "info" I don't really know if it's necessary or not. I didn't want to check, so I just include it.

<p>The function <i>get_root_info</i> is not very interesting. I won't show it here.

<p>The function <i>get_files</i> is funny. It looks for all the <i>.obj</i> files in the current directory and puts them in a list. Then it sorts the list alphabetically, to get the object in the correct order for later on. Here it comes.

<p>

<pre>

```
void get_files(pdf_file *f)
{
    DIR *d = opendir(".");
    struct dirent *e;
    struct stat s;
    obj *o;
    if (!d) ERR;
    while (1) {
        e = readdir(d); if (!e) { if (errno) ERR; break; }
        if (e-&gt;d_name[0] &lt; '0' || e-&gt;d_name[0] &gt; '9' || !strstr(e-
&gt;d_name, ".obj"))
            continue;
        o = new_obj(f);
        o-&gt;file[FILENAME_MAX-1] = 0;
        strncpy(o-&gt;file, e-&gt;d_name, FILENAME_MAX-1);
        if (stat(e-&gt;d_name, &s)) ERR;
        o-&gt;size = s.st_size;
    }
    closedir(d);
    qsort(f-&gt;o.o, f-&gt;o.size, sizeof(obj), cmp);
}
</pre>
```

<p>See the use of <i>qsort</i> to, well, sort the list of object. The <i>libc</i> is friendly, sometimes.

<p>And we finish that with <i>dump</i> which processes objects in order. It stores the number of bytes written so far in the <i>s</i> variable, and updates <i>o->offset</i> for each object so that the trailer gets a correct value for the starting offset of each object.

<p>

<pre>

```
void dump(pdf_file *f)
{
    FILE *fi;
    size_t s = 0;
    obj_list *l = &f-&gt;o;
    obj *o;
    int i;
    int c;
    s = printf("%PDF-1.4\n%\x88\x89\x90\x91\n");
    for (i = 0; i &lt; l-&gt;size; i++) {
        o = &l-&gt;o[i];
        o-&gt;offset = s;
        fi = fopen(o-&gt;file, "r"); if (!fi) ERR;
    }
}
```

```

while (1) {
    c = getc(fi); if (c == EOF) { if (errno) ERR; break; }
    if (fwrite(&c, 1, 1, stdout) != 1) ERR;
    s++;
}
if (fclose(fi)) ERR;
}
printf("xref\n0 %d\n0000000000 65535 f \n", l->size + 1);
for (i = 0; i < l->size; i++) printf("%10.10d 00000 n \n", l->o[i].offset);
printf("trailer\n");
printf("<</Size %d/Root %d 0 R/Info %d 0
R&&\nstartxref\n%d\n%%EOF\n",
    l->size + 2, f->root, f->info, s);
fflush(stdout);
}
</pre>

```

So this program is the most complex one, indeed. We need to store a list of objects. You just can't process data on the fly. The trailer forces us to "remember" some stuff for later use.

See the structures I had to create.

```

<p>
<pre>
typedef struct {
    char file[FILENAME_MAX];
    size_t size;
    size_t offset;
} obj;

typedef struct {
    obj *o;
    int size;
    int maxsize;
} obj_list;

typedef struct {
    int root;
    int info;
    obj_list o;
} pdf_file;
</pre>

```

I wonder if a program using recursion rather than a list would not be possible. Would it be smaller? more efficient? more elegant? Well, this is not my style of programming, so I won't dig too much in that direction. People more in the functional programming may have some funny ideas there.

And how to pilot this program?

Easy.

```
<pre>
cd out; ../glue > TTT
</pre>
```

<p>And there you are. A nice, anonymous PDF file called <i>TTT</i>. Why this strange name? Well, once again, laziness. I had to find one, I picked up something in my brain, for no particular reason. Habits, maybe. Weird habits? Maybe.

<h2>Automatizing the process</h2>

<p>Obviously each time something new appears in there I had to run the process from scratch once again.

<p>There were many bugs. Some at the syntactic level, like calling a wrong function, wrong arguments, this kind of things. Some others were due to a bad understanding of the PDF specification. Some others to bad shell's knowledge.

<p>There are many bugs that pop up here and there. A hacker lives with those beasts, looks for them, eliminates them. You get used to it. Whatever procedure you put in yourself, you won't avoid bugs. So you need to develop here again a good culture of potential problems.

<p>One nasty source of problems comes from the tools you use and the error messages they throw at your face. I specially think about the compiler. Sometimes you just forget a <i>;</i> and all in a sudden you get hundreds and hundreds of unrelated errors. Or you forget a <i>}</i>. Or some minor stuff like that.

<p>There again you live with that. You compile often, so you are sure a compilation error can only originate from the few lines you introduced and not somewhere else.

<p>And you debug harder bugs too. It may be frustrating sometimes. Getting good habits, a good culture, takes time and devotion. You have to bypass frustration and just keep going, up to success. I mean, this is not research. There is a solution, you know that from the beginning. So it's just a matter of doing things properly to get the expected result. That is what is so cool with hacking. The results are there, necessarily.

<p>And so, to get back to automation, I could have written a "pilot" in C. But, traditionally, this is done in shell script.

<p>We already saw the various shell's commands, but to recapitulate, here comes <i>transform.sh</i>.

```
<p>
<pre>
rm -f out/*
cp trailer.obj out
```

```

for i in ????.obj; do ./deflate < $i > out/$i; done
(cd out; find -empty -exec cp ../{} . \;)
(cd out; for i in `grep -i limited *|cut -d : -f 1`; do ../empty $i > $i;
done)
#(cd out; for i in `grep -i limited *|cut -d : -f 1`; do cp $i $i.emp; done)
for i in *.n; do cp $i out/`echo $i|sed -e "s/\.n//"`; done
(cd out; for i in ????.obj; do ../inflate $i; done)
(cd out; ../glue > TTT)
</pre>

```

<p>Each line of this file would require several dozens of C code. That's why shell programming is something to learn too. It speeds things up for basic data manipulation.

<p>But some treatments are beyond the shell and various command line's programs' ability and you need to write more complex program.

<p>I do it in C, because that's what I know best. Several years of practice shaped my brain to easily think in C.

<p>The main problem you face when you hack is bugs. And knowing a language, practicing it over the years, gives you familiarity. You know where to look for bugs. You know that when you looked here and found nothing then there is nothing wrong there. This knowledge you don't get for systems and languages you don't know. So you can't narrow the source of a problem. You can't become <i class="white">certain</i>.

<p>Well, that's it.

<h2>Conclusion</h2>

<p>Do these programs "think?"

<p>They manipulate data, extract information. As you and I do through our sensory system. Is it comparable?

<p>Some might say they are not complex enough. They just apply a given procedure.

<p>How are you certain you don't work this way, huh?

<p>I don't pretend the programs here think. It's just not that simple to give a convincing <i class="white">no</i>. Take that question as an opportunity to think about yourself, the world, life, the meaning (if any) in there and so on.

<p>Or not. I mean, it could be seen as some kind of madness or sick attitude or whatever. It's just a tough experiment, based on "reality" (as long as a hack is "real").

<p>Anyway, I'm done with this hack. Let's move on to something else. (The magic of hacking. It never ends.)

<p>(By the way, what is IEEE 754? Short answer: dig the web. Longer answer: it's just a standard to represent floating point numbers in computers. And yes, normal people would think it's very boring. But I am not normal and

neither are you! Are you?)

<hr>

<p>Contact: sed@free.fr

<p>Created:

Mon, 06 Jul 2009 17:03:23 +0200

Last update:

Mon, 06 Jul 2009 17:03:23 +0200

</body>

</html>